

CS7IS3: Assignment 1 Crandfiled search engine

Mingwei Shi
Trinity college dublin
mshi@tcd.ie

I. INTRODCUTION

Parser-Index-Search -Evaluation

II. IMPLEMENTATION

A. Parser

Finite state machine to process

B. Indexing

Feature type index:

Only two feature are indexed by Textfile that could be tokenized including title and description. Since the rest of feature should not be separate

Customiser :tokenize+lowercase+possession removal+ large stop words removal.

C. Searcher

Similarity use: VSM;BM25;LMDIRCHATE

III. EVALTION

In this section, I would like to identify the discriminating features combination that has an impact on mean precision value, what is the influential analyzer that has an impact on mean precision value, and what is the influential similarity matching algorithm that has impact on the mean precision value.

a1 for analyzer : Standard

a2 for analyzer : Customized

s1 for Similarity : Vector space

s2 for Similarity : BM25

s3 for Similarity : LMDirchite

f2 for Description feature only(.W) : Vector space

f3 for combining with title+ description(.T and .W) feature only(.W) :

Observation for analyzer: When add more stopword,the recall decrease as the key term are smaller;

Observation for features:The content of description(.W) is so informative that introducing extra content such as title(.T) would not add the precision.As the information entropy might not increase.

Observation for similarity: BM25 is the best

REFERENCES

TABLE I
COMPARISON BETWEEN DIFFERENT ANALYZERS, SIMILARITY, AND
FEATURE.

Acronym Name	MAP
a1-s1-f2	0,2831
a1-s1-f3	0,2831
a1-s2-f2	0,2945
a1-s2-f3	0,2945
a1-s3-f2	0,231
a1-s3-f3	0,231
a2-s1-f2	0,1667
a2-s1-f3	0,1667
a2-s2-f2	0,1616
a2-s2-f3	0,1616
a2-s3-f2	0,1121
a2-s3-f3	0,1121